

AUTOMATIC OBJECT MIGRATION ALTERNATIVES FOR AD HOC AND SENSOR NETWORKS

Rimon Barr T. W. Danny Kim John C. Bicket Emin Gün Sirer

Computer Science Department
Cornell University
Ithaca, NY 14853
{barr, egs}@cs.cornell.edu; {tk65, jcb35}@cornell.edu

November 4, 2001

Developing applications for ad hoc and sensor networks poses significant challenges. Many interesting applications in these domains entail collaboration between components distributed throughout an ad hoc network. Defining these components, optimally placing them on nodes in the ad hoc network and relocating them in response to changes is a fundamental problem faced by such applications. Manual approaches to object migration are not only platform-dependent and error-prone, but also needlessly complicate application development. Further, locally optimal decisions made by applications that share the same network can lead to globally unstable and energy inefficient behavior.

In this paper we describe the design and implementation of a distributed operating system for ad hoc and sensor networks whose goal is to enable power-aware, adaptive, and easy-to-develop ad hoc networking applications. Our system achieves this goal by providing a single system image of a unified Java virtual machine to applications over an ad hoc collection of heterogeneous nodes. It automatically and transparently partitions applications into components and dynamically finds a placement of these components on nodes within the ad hoc network to reduce energy consumption and increase system longevity. This paper outlines the design of our system and evaluates two practical, power-aware, online algorithms for object placement that form the core of our system. We demonstrate that our algorithms can increase system longevity by a factor of four to five by effectively distributing energy consumption, and are suitable for use in an energy efficient operating system in which applications are distributed automatically and transparently.

Keywords: ad hoc networks, power-aware object migration, single system image, application partitioning, operating system

1 Introduction

Ad hoc networks, formed through the collaboration of intelligent mobile nodes over wireless links, simultaneously promise a radically new class of applications and pose significant challenges for application development. Recent advances in low-power, high-performance processors and medium to high-speed wireless networking have enabled new applications for ad hoc and sensor networks, ranging from large-scale environmental data collection to coordinated battlefield and disaster-relief operations. Many interesting applications entail collaboration between components distributed throughout an ad hoc network. For example, sensor networks are often composed of three types of components: sensors acting as data sources, information consumers operating as data sinks, and numerous filters in-between for performing application-specific data processing. While the data sources and sinks may be coupled tightly to the nodes to which they are attached, there is often a high degree of freedom in the placement of the data processing components. This freedom, coupled with the dynamic environment posed by ad hoc networks, both enables adaptive applications and makes it difficult to find the optimal distribution of application components among nodes.

Adapting to dynamically changing conditions by changing the distribution of functionality across a network is critical for many distributed ad hoc networking applications [Satyanarayanan 96]. For example, the resources available to components at each ad hoc node, in particular the available power and bandwidth, may change over time and necessitate the relocation of application components. Further, event sources that are being sensed in the external environment, such as tracked objects or chemical concentrations, may move rapidly, thereby shift network loads and require applications to adapt by migrating components. Finally, application behavior might change, as in the transition from defensive to offensive mode in a battlefield application, modifying its communication pattern and necessitating a reorganization of its deployed components within the network.

Currently, ad hoc networking applications either rely on a static assignment of components to nodes or use *ad hoc*, manual policies and mechanisms for migrating objects in response to change. A static assignment of functionality to nodes simplifies application design by obviating migration and reduces meta-traffic in the network by eliminating object mobility, but it also leads to non-adaptive, fragile and energy-inefficient systems. The overall application will stall as soon as the critical nodes on the dataflow path run out of power or move out of transmission range. Manual approaches to object mobility suffer from being hard to develop, error prone and platform specific. Each application using this approach needs to re-implement the same migration, monitoring and communication mechanisms, correctly, on every platform. Further, locally optimal policies pursued by individual applications may lead to globally unstable and energy-inefficient behavior when multiple applications share and interact on the same ad hoc

network. In essence, this approach suffers from building on an abstraction-level that is too low. A high-level operating system that provides the requisite mechanisms and policies for code mobility would not only simplify application development, but also ensure the integrity of system-wide goals in the face of multiple applications competing for resources.

In this paper, we outline the design of a single system image operating system for ad hoc networks. The goals of our system are to enable the construction of applications with the following properties:

- **Adaptive:** Applications should, with minimal effort, be able to respond to changes in their environment, their communication pattern, and the availability of resources in the network.
- **General purpose:** Applications should be able to run on both ad hoc and fixed networks. Porting an existing monolithic application to execute efficiently on an ad hoc network should require little effort.
- **Platform independent:** Applications should be able to execute on ad hoc networks of nodes with heterogeneous resources and capabilities.
- **Efficient:** Policies and mechanisms used for adaptation in the systems layer should not require excessive communication or power consumption. The default policies and mechanisms should yield good power utilization and maximize total system lifetime.

Our operating system meets these goals by providing the illusion of a single, unified Java virtual machine over an ad hoc network of heterogeneous, physically separate, mobile hosts. Our system consists of a static application partitioning service that resides on border hosts capable of injecting new code into the network, and a runtime on each node that performs dynamic monitoring and object migration. The static partitioning service takes regular Java applications and converts them into distributed components that communicate via RMI by rewriting them at the bytecode level (Figure 1). The code injector then finds a suitable initial layout of these objects and starts the execution of the application. The runtime monitors the performance of the application and migrates application objects when doing so would benefit the system.

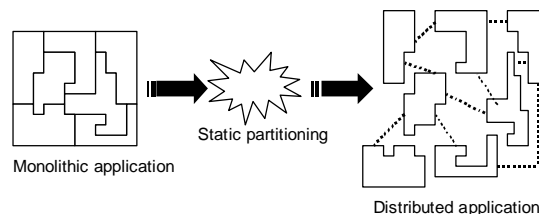


Figure 1: A static partitioning service converts monolithic Java applications into distributed applications that can run on an ad hoc network and transparently communicate via RMI.

The core of our system consists of algorithms for deciding when and where to move application components. While our system is designed such that these algorithms can be transparently replaced to optimize for differing goals, such as minimizing application latency, response time, or bandwidth consumption, in this paper we tackle what we believe to be the most important goal in energy-constrained ad hoc networks of mobile hosts. Namely, we examine how to maximize total application lifetime by

utilizing energy more efficiently. We present two practical, online algorithms, named *NetPull* and *NetCenter*, for finding a distribution of application components on nodes in an ad hoc network that increases total system lifetime by increasing energy utilization (Figure 2). We evaluate these algorithms in the context of a generic sensing application and examine their impact on *system longevity*, which we define as the length of time that a generic sensing application can maintain sensor coverage above a given threshold area. Both algorithms operate by dividing time into epochs, monitoring the communication pattern of the application components within each epoch, and migrating components at the end of the epoch when doing so would result in more efficient power utilization. We show that *NetCenter* or *NetPull* achieve a factor of four to five improvement in system longevity over naïve or static partitioning techniques.

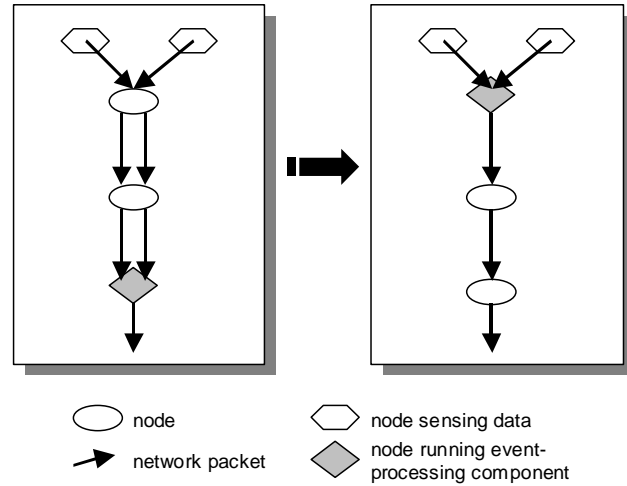


Figure 2: Automatically migrating components closer to their data sources in a sensor network increases system longevity and decreases power consumption by reducing total network communication cost.

This paper makes three contributions. It outlines the design and implementation of a single system image operating system for ad hoc networks, where the entire network appears to be a large Java virtual machine to applications, whose components are partitioned among the nodes automatically and migrated transparently. Secondly, we propose two practical, adaptive, online algorithms for deciding when and where to move application components. Finally, we demonstrate that these algorithms achieve high-energy utilization, extract low overhead, and improve system longevity, and are thus suitable for use in an operating system for transparent, automatic object migration.

In the next section, we describe related work on operating system support for ad hoc networks and their applications. Section 3 outlines our system implementation, including the code partitioning service and runtime support for object migration. Section 4 describes *NetPull* and *NetCenter*, two algorithms for automatic object migration. Section 5 presents our evaluation of these algorithms and of the automatic application partitioning. Section 6 describes our plans for future work. We summarize our contributions in Section 7.

2 Related Work

Prior research on ad hoc networks has concentrated on ad hoc routing algorithms such as DSDV [Perkins & Bhagwat 94], DSR [Broch et al. 96], AODV [Perkins 97], ZRP [Haas & Pearlman 98], TORA [Park & Corson 98], GeoTORA [Ko & Vaidya 00], PARO [Gomez et al. 01], OLSR [Jacquet et al. 98], LANMAR [Pei et al. 00], AMRIS [Wu & Tay 99], CEDAR [Sivakumar et al. 99] and others [MANET 01, Royer & Toh 99]. These routing algorithms move data from a source to a given destination(s) as efficiently as possible. They assume that the application was structured *a priori* by the programmer to operate optimally in an ad hoc environment and that the communication pattern of the application is fixed. Specifically, the route is the focus of the optimization, not the source and destination endpoints. Our system complements routing algorithms by moving application code around the network, as well as data. Bringing application components closer to the data sources radically alters the overall communication pattern of the application, and can increase system and application longevity.

Others have experimented with single system image (SSI) operating systems for wired, tightly coupled networks of workstations. Sprite [Ousterhout 88], V [Cheriton 88], Ameoba [Steketee 95], Accent [Zayas 87], and LOCUS [Popek & Walker 85] implement native operating system facilities for migrating processes between nodes on a tightly coupled cluster. Most recently, the cJVM [Aridor et al. 99] and JESSICA [Ma et al. 99] projects provide a similar Java virtual machine-based, single system image to their applications. Others, including Condor [Litzkow et al. 97], libckpt [Plank et al. 95] and CoCheck [Stellner 96], provide user-level mechanisms for process migration without operating system support. These projects target high-performance, well-connected clusters, and face an entirely different set of challenges; the main goals of these SSI systems were load balancing and performance in a local area network for interactive desktop programs or compute-intensive batch jobs. In contrast, our system targets wireless multi-hop networks, where utilizing power effectively and maximizing system longevity is more important than traditional application performance.

Some recent systems have examined how to spatially partition applications within a wired network. The Coign system [Hunt & Scott 99] has examined how to partition COM applications between two tightly interconnected hosts within a local-area network. Coign performs static spatial partitioning of office applications via a two-way minimum cut based on summary application profiles collected on previous runs. Extending this work, the ABACUS system [Amiri 00] has shown that incorporating dynamic information into component placement decisions can further improve application performance. Our work shares the same insight as Coign, in that we are also interested in the automatic partitioning and relocation of application components, but differs in that it partitions applications across any number of nodes, instead of just two, and moves application components dynamically, in response to changes in the

network instead of computing a static partitioning from a profile. We further differ from both these systems with our optimizations targeted at power conservation and system longevity.

Various mobile computing projects have looked at how to structure mobile applications. The Emerald system [Jul et al. 88] provides transparent code migration for code components written in the Emerald language, where code migration is directed by source-level programmer annotations in Emerald. The xDU project [Gehani 97] provides a similar framework for explicitly partitioning applications into distributable units using code annotations in Java. Legion [Lewis & Grimshaw 95] is a language independent, scalable, object-oriented operating system for wide-area infrastructure networks. They differ fundamentally from our system in that they do not provide an encompassing operating system platform, require explicit programmer control to trigger code migration, do not support an ad hoc network model and target traditional applications.

Other mobile computing projects have focused on creating mobility toolkits to facilitate the construction of mobile applications. The Rover toolkit [Joseph et al. 95] provides relocation and messaging services for disconnected and mobile operation. More recent work includes the creation of toolkits specifically for ad hoc environments and applications. The Mobeware toolkit [Campbell 98] provides an adaptive-QoS programming interface, which assists in the management of object flows. XMIDDLE [Mascolo 01] assists with data management and synchronization. Our research takes a systems approach instead of a programmer driven toolkit. Our system automatically manages the shared network and energy resources among ad hoc sensor applications. A benefit of this approach is that it obviates the need for applications to be recoded against a specific toolkit interface, since applications are executed without modification atop our operating system. Furthermore, this approach ensures that applications, regardless of which toolkits they use, behave in a cooperative manner.

Active networks [Tennenhouse & Wetherall 96] are a related technique for distributing computation within a network on a per-packet or per-flow basis. Our system differs from active networking efforts in that it offers higher-level, encompassing operating system services for ad hoc applications. Its resource accounting is based on computational nodes and application components instead of packets or flows. It is, however, complementary to active networks, in that it may use an active network as a low-level mechanism for code distribution. In a similar vein, recent work on directed diffusion [Heidemann et al. 01] argues for the labeling of data, then routing and filtering packets in the network using stateless components specified in a constrained language. In contrast, we provide a complete system that automatically migrates components in the network, and does so in an energy-conserving manner. Furthermore, our components may be stateful and can therefore perform aggregation and more

sophisticated operations in addition to basic routing and packet filtering. Our system supports the automatic migration of this state.

Some mobile applications can naturally be structured as instances of distributed leader election. Much prior work has investigated the problem of selecting a leader node among a set of replicas [Chang and Roberts 79, Garcia-Molina 82, Awerbuch 87, Ostrovsky 94, Sayeed et al. 95, Brunekreef et al. 96, Singh 96, Russell et al. 99], including recent work in an ad hoc network setting [Malpani et al. 00]. However, these algorithms do not account for the cost of keeping the state of object replicas consistent, nor the effect of the messages exchanged for leader election on power consumption and system longevity. For applications with stateful components these costs and effects can be substantial.

Lastly, prior work has examined how to minimize power consumption within an independent host by manipulating the processor clock speed [Grunwald et al. 00, Weiser et al. 94]. Our system is complementary to this work and opens up further opportunities for minimizing power consumption by shipping computation out of hosts limited in power to less critical nodes.

3 System Design and Implementation

In this section, we outline the design of our system and discuss some implementation decisions. We focus on those aspects related to object mobility, specifically the manner in which applications are partitioned, the runtime support for migration and the object affinity interface, which explicitly restricts automated migration decisions.

3.1 Application Partitioning

Our system provides an automatic code partitioning service that divides regular, monolithic Java applications into components that can migrate and execute over a network. Our approach to partitioning applications statically is patterned after the Distributed Virtual Machine paradigm [Sirer et al. 99], and provides three advantages. First, the complex partitioning services need only be supported at code-injection points, and can even be performed offline. Second, since the system operation and integrity do not depend on the partitioning technique, users can manually partition their applications into arbitrary components, if they so choose. Finally, this technique provides a convenient, default mechanism for transitioning legacy, monolithic applications to execute over ad hoc networks.

Partitioning applications involves separate modifications for object creation, invocation and field access, respectively. Object creation instructions are replaced by calls to the local runtime, which is responsible for selecting an appropriate node and creating a new instance of the object on that node. This operation returns a handle to the remote object, which can then be used for subsequent object invocations.

Object invocations are modified to go through an RPC mechanism [Birrell & Nelson 84] to invoke remote objects. Field member accesses are rewritten as invocations of synthetic accessor methods. We use the Java RMI interface for remote object invocation [Harold 00]. The division is performed statically along object interfaces at class granularity¹. Thus, the transformation preserves class interfaces and retains type safety. It works at the bytecode level without requiring source access.

3.2 Runtime Support for Object Migration

Our system provides a small runtime to manage individual nodes. The runtime facilitates communication between application components. In addition, it transparently monitors application behavior to enable intelligent migration decisions. It is invoked in response to three events: object creation, invocation and migration.

In order to create a new instance of an object, an application contacts the local runtime and provides the type of the object to create. The runtime then has the option of placing the newly created object at a suitable location with little cost. It may choose to locate the object on the local node, at a well-known node or at its best guess of an optimal location within the network. In our current implementation, all new objects are created at the local node. We chose this approach for its simplicity, and rely on our dynamic object migration algorithms to find the optimal placement of objects over time. The application binaries, containing all of the object code, are distributed to all nodes at the time that the application is introduced into the network. For large networks, or when code changes dynamically, this code-loading can also be performed in a lazy manner on demand at each node. Once created, the runtime simply initializes the object by calling its constructor and returns a remote handle of the object to the caller.

The runtime handles object invocations by locating the appropriate stubs, marshaling arguments, and performing the right RPC invocation via RMI. Our runtime independently keeps track of object locations in the OS layer. But we assume the presence of a standard ad hoc routing protocol, like AODV or DSR, below our runtime to provide message routing. When an object migrates, it sends a multicast message to all of the components that have outstanding handles for itself and informs them of its new IP address. This multicast is tightly integrated with the underlying routing protocol to inform remote stubs of the new location of the object while simultaneously establishing a live route. In the case of AODV, the multicast packet contains a piggybacked RouteReply and causes intermediate nodes to acquire a distance metric and route for the new destination node. [*Note to reviewers: Our runtime implementation does not currently perform this optimization; however, we simulate its effects in the evaluation section.*]

¹ Consequently, application components correspond to Java objects in our system, and we use the two terms interchangeably throughout this paper.

Finally, the runtime has a background thread that listens to incoming migration requests, receives new objects over the network and instantiates them. The class bytecode for each object is usually aggressively injected into the network and exists at each node. If the binaries do not reside at the node when instantiation is requested, they are retrieved by the class loader.

3.3 Object Affinity

In order to gain acceptance, any automatic scheme for code migration should be at least as good as any manual scheme for migrating components. We ensure that our system is no worse than manual schemes by providing an explicit interface by which application writers can manually direct the object placement performed by our system. The interface enables a programmer to explicitly establish affinities between components and ad hoc nodes.

We provide two levels of affinity. Specifying a *strong* affinity between a component and a node effectively anchors the code to that node. This is intended for attaching components like device drivers to the nodes with the installed device in them. Specifying a *weak* affinity immediately migrates the component to the named node, and allows the automated code placement techniques described herein to adapt to the application's communication pattern from the new starting point. Note that today's manually constructed applications correspond to the use of strong affinity in our system, wherein components are bound to nodes unless explicitly moved. The result of overusing strong affinity is a fragile system, where unforeseen communication and mobility patterns can leave an application stranded. While we provide these primitives, we do not advocate their use and believe that automated techniques can outperform manual efforts to place components.

3.4 Implementation Details

We chose to implement our system in Java. Among other benefits, the Java platform allows us to rewrite binaries without source access and facilitates object migration as described above. Deploying applications in ad hoc networks requires operating within the resource limitations of small, cheap sensors. Yet, Java virtual machines on the desktop can have excessive resource requirements, necessitating a fast processor for just-in-time compilation, RAM for expanded object storage, and persistent storage for the system libraries. Recent work [Surer et al. 99] has directly addressed these issues and proposed a new system architecture for virtual machines that can reduce all three of these resource requirements. Indeed, there are Java Card virtual machines in existence today that are based on a similar, partitioned service architecture that fit on a flexible credit card and cost a few dollars. We assume that running a Java interpreter, or equivalent functionality, on sensors is a solvable problem and concentrate on the distributed coordination of applications.

Our current network stack, including a reliable socket layer and a full AODV implementation, is implemented in user space on top of UDP datagrams. We are currently implementing this stack inside the Linux kernel to improve performance. In this implementation all packets sent via standard sockets are processed by an AODV layer. A user daemon implements the AODV protocol, and intercepts ARP packets on the interface using route requests and route reply packets to populate the IP route table instead.

4 NetPull and NetCenter

In the preceding section, we described how our system provides object mobility. In this section, we describe two algorithms, named *NetPull* and *NetCenter*, which use these mechanisms to increase system longevity.

Both *NetPull* and *NetCenter* share the same basic approach. They shorten the mean path length of data packets by automatically moving communicating objects onto nodes that are topologically closer together. They perform this by profiling the communication pattern of each application in discrete time units, called epochs. In each epoch, every runtime keeps track of the number of incoming and outgoing packets for every object. At the end of each epoch, the migration algorithm decides whether to move that object, based on its recent pattern of behavior. Under both algorithms, the decision is made locally, based on information collected during recent epochs at that node. The epochs need not be synchronized across nodes, as the decisions to move objects are made locally and independently. *NetPull* and *NetCenter* differ in the type of information they collect and how they pick the destination host. Depending on the environment, one may be easier to implement.

NetPull collects information about the communication pattern of the application at the physical link level, and migrates components over physical links one hop at a time. This requires very little support from the network; namely, the runtime needs to be able to examine the link level packet headers to determine the last or next hop for incoming and outgoing packets, respectively. For every object, we keep a count of the messages sent to and from each neighboring node, identified by its network interface MAC address. Each packet from a neighboring host is assigned a net force of one unit in the direction of that host. At the end of an epoch, the runtime examines all of the forces on the object and moves it one hop in the direction of greatest pull.

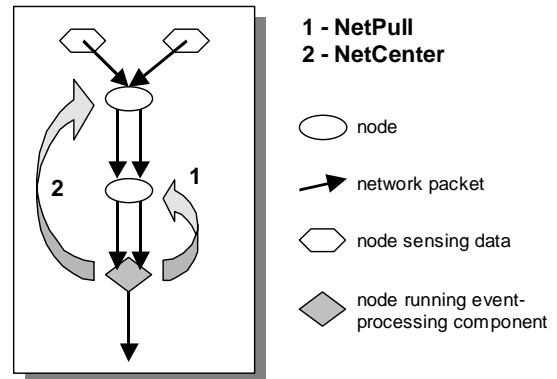


Figure 3: NetPull (1) moves one hop towards the source of data whereas NetCenter (2) moves directly to the source of most packets.

NetCenter operates at the network level, and migrates components multiple hops at a time. In each epoch, *NetCenter* examines the network source addresses of all incoming invocations, and the destination addresses of outgoing invocations, for each object. The network source and destination addresses are often part of the packets sent during the RMI operation. For instance, they are often part of the IP source and destination fields, and thus can be collected without placing an extra burden on the network. At the end of an epoch, *NetCenter* finds the host with which a given object communicates the most and migrates the object directly to that host, skipping over any intermediate nodes.

Both of these algorithms improve system longevity by using the available power within the network more effectively. By migrating communicating components closer to each other, they reduce the total distance packets travel, and thereby reduce the overall power consumption. Further, moving application components from node to node helps avoid hot spots and balances out the load on the network. As a result, both algorithms significantly improve the total system longevity for an energy-constrained ad hoc network.

5 Evaluation

In this section, we focus on the evaluation of the key mechanisms of our operating system. We show that *NetPull* and *NetCenter* achieve good energy utilization, improve system longevity, and are thus suitable for use in a general-purpose, automatic, transparent object migration system. We describe our networking model, benchmark application, and simulation framework, and then discuss the results in detail. We also evaluate the performance of automatically partitioned code against code that is manually generated, and observe that the code transformation overhead is insignificant.

5.1 Network Model and Benchmark Application

Our system targets general-purpose ad hoc networks. These networks consist of geographically distributed nodes communicating over wireless links. All nodes have the same communication radius and they are connected to the fixed networking infrastructure via a single, centrally placed node. Each node initially stores a fixed, finite amount of energy. Sending a packet between any neighboring nodes exacts a constant communication cost. We currently assume that the cost of local computation on a host is negligible in comparison with packet transmission costs.

We examine a generic, reconfigurable sensing benchmark we developed named SenseNet. This application consists of sensors, condensers and displays. Sensing components are fixed at particular ad hoc nodes, where they monitor events within their sensing radius and send a packet to a condenser in response to an event. Events on the field are generated with a random waypoint model, and each

condenser is responsible for a subset of nodes in a region of the field. Condensers can reside on any node, where they process and aggregate sensor events, as well as filter noise. In fact, condensers start at the central node, and migrate towards their event sources as the application runs. The display runs on the central node, extracts high-level data out of the sensor network and sends it to the wired network. We have arbitrarily defined failure for this application as the point when half of the field is no longer being sensed, that is, only half of the field area is within the sensing radius of at least one live sensor which can communicate along some functioning route with the central node.

5.2 Simulation Framework

We developed a fast, scalable, packet-level, statically parameterizable ad hoc networking simulator in order to simulate large networks. In the style of ns2 [McCanne & Floyd 01], the simulator accounts for all communication costs, including DSR routing and route repair overhead. Unlike ns2, all setup, configuration and customization of our simulator is done at compile-time instead of run-time. This highly customized, scenario-specific simulation strategy allows us to simulate large-scale sensor networks. Below, we present results for a network of 3600 nodes.

The simulator models the movement of every unicast and multicast packet and properly incurs the cost of moving a condenser and notifying all its sensors of the new location. We initialize the simulator with a uniform distribution of nodes on a plane, and vary parameters such as noise levels, field size, density, battery power, and communication and sensing radii. Sensing events are generated at random locations on the field, and with random durations and velocity vectors. Sensors that are in range detect these signals and generate application events. They can also, with small probability, generate fictitious events due to sensor noise.

5.3 Algorithms

We compare four different algorithms for transparent object migration:

- **Static** corresponds to a static, fixed assignment of objects to nodes within the network. Our components remain at the home node for the entire duration of the simulation.
- **Random** is designed to overcome a fundamental shortcoming of static placement, namely that the entire system can no longer function once key nodes are unreachable, by moving components around at random.
- **NetPull** looks at link-level data from its neighbors and moves to the most active neighbor.
- **NetCenter** keeps tracks of message sources and moves directly to the node with greatest activity.

5.4 Simulation Parameters

A simulation of a complex system such as this requires many parameters. We summarize them here. In the following experiments, we examine a simulated network on a field of 300 by 300 distance units. Total number of nodes is 3600, corresponding to a density of 0.04. Node sensing radius is 20 units; communication radius is 10. Sensors generate spurious messages due to noise with probability 1% in each epoch. Sending a packet incurs unit energy cost per hop, and the initial battery capacity is arbitrarily chosen to be 1000. We examined, but do not present results from, simulations with the following parameters: density={0.02, 0.03, 0.04, 0.05, 0.06}, noise level={0.01, 0.05, 0.10}, initial battery power={1000, 2000}. The choice of a particular density, noise level, or initial battery power does not materially impact our results. The choice of epoch duration is arbitrary and application-dependent from the point of view of the system and should be made in consultation with applications. Each epoch contains at least one event in our simulations. An event may span up to 10 epochs and move through the field with a uniformly chosen velocity between 0.0 and 2.0 distance units per epoch. We assume that nodes are stationary after the initial deployment, though none of the nodes make any assumptions about geographical location of other nodes. Every data point represents an average of five runs.

5.5 Results and Discussion

Figure 4 illustrates the impact of our algorithms on system longevity. The graph shows the number of epochs before application failure for each of the migration strategies. For this and subsequent graphs, we defined the application failure point to be when less than 50% of the field can be sensed, either because sensors have been completely drained of energy or have been disconnected from the central node by other drained nodes. We ran the simulation with different application failure thresholds, with similar results: *NetPull* and *NetCenter* lengthen the operational lifetime of the system by a factor of four to five. By taking application communication behavior into account, these policies migrate objects to reduce the

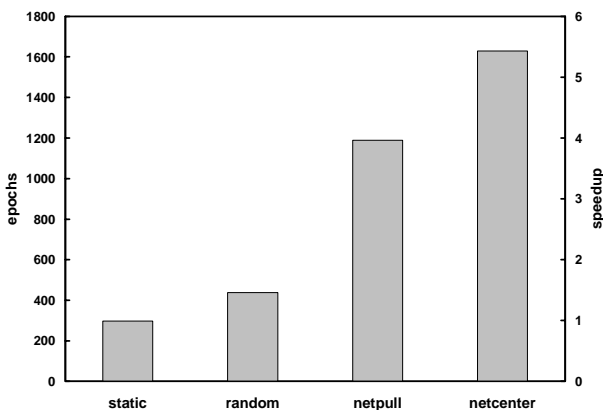


Figure 4: System longevity improvement

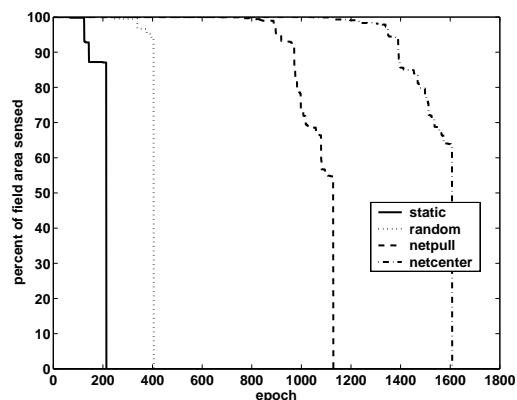


Figure 5: Sensor coverage degradation over time

mean path length of packets, thereby reducing energy consumption and extending system longevity.

Figure 5 shows how quickly the field coverage degrades when components are assigned to nodes in a manner that is oblivious to the underlying application communication pattern. In contrast, *NetPull* and *NetCenter* migrate components close to the source of events, thus preserving connectivity by shedding communication load from the critical nodes in the network. Furthermore, we see a sharp drop in sensor coverage from around 90% coverage to zero coverage for *Random* and *Static*. This happens when the last live node in the ring around the critical central node is drained, disconnecting the entire network that was previously communicating with the central node using this node as the last hop. *NetPull* and *NetCenter* exhibit a more gradual loss of connectivity. This early warning of impending failure is a useful system property, allowing nodes to be replenished or replaced.

Figure 6 examines the number of nodes that run out of power as a function of time. The curve ends when the ad hoc sensor network application can no longer cover more than the threshold area. The endpoints of the different curves demonstrate that static and random object placement schemes fail sooner, as a result of fewer failures in the network. In contrast to the fragility of static and random object placement, informed active migration algorithms can continue to function even in the presence of large numbers of failures in the network. This is because active migration algorithms avoid creating communication hotspots around critical nodes, maintaining connectivity in the network for longer. Both *Static* and *Random* reach failure prematurely: few nodes are drained, yet sensor coverage has been eroded, as is also evident in Figure 5. In contrast, *NetPull* and *NetCenter* drain fewer nodes and preserve connectivity.

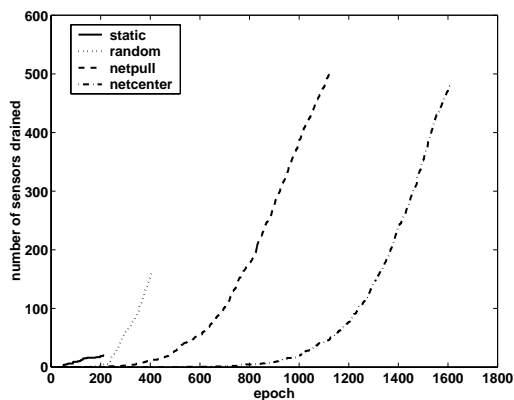


Figure 6: Sensor drainage over time

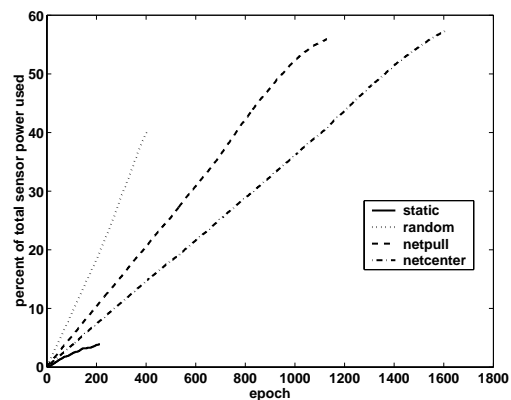


Figure 7: Field energy used over time

Figure 7 shows a counter-intuitive result: *NetPull* and *NetCenter* lie between *Static* and *Random* in terms of energy consumption per epoch. The low slope for *Static* shows that this algorithm uses little energy per epoch because it does not expend any power on active object migration. Most of this power usage is concentrated in a ring around critical nodes, however, and the application terminates early,

leaving more than 95% of the total energy on the field uninitialized. *Random* avoids some communication hotspots by simply moving objects around randomly, and outlasts *Static*, as seen in Figure 5. However, *Static* and *Random* both fail early. *NetPull* and *NetCenter* also consume energy to move components, but consume less energy and last longer than *Random*, because they make informed migration decisions that reduce the mean path length of communication. Overall, they outlast both *Random* and *Static* showing that it is worthwhile to spend energy to move components, when taking application communication behavior into account.

The graph of disconnected nodes over time, shown in Figure 8, indicates that the number of disconnected nodes increases more gradually for *NetPull* and *NetCenter* because they distribute load more evenly across the network and avoid completely draining key nodes. In the case of *Static* and *Random*, even though only a small number of nodes are drained, they are all located around the home node and thus quickly disconnect the entire field from the wired network.

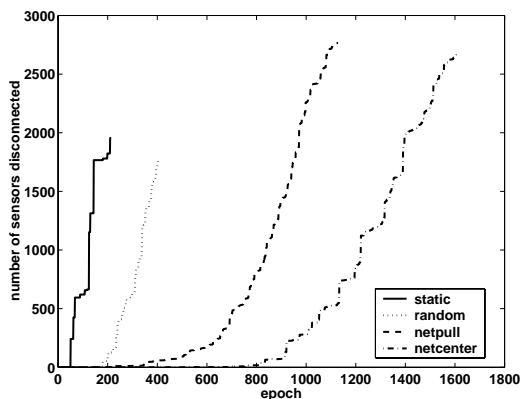


Figure 8: Sensors disconnected over time

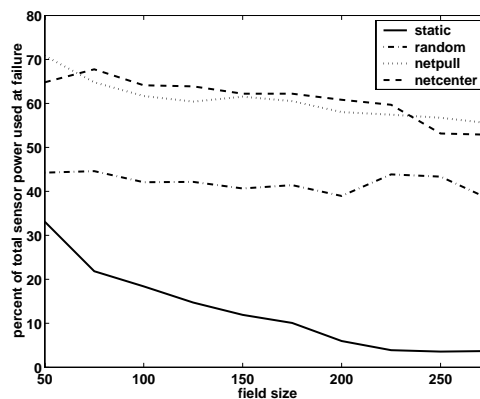


Figure 9: Energy used at breakdown versus field size

Figure 9 plots the total energy remaining on the field at system failure time versus the field size. All algorithms, except for the static placement of components, are unaffected by variations of field size, and will scale to large networks. *Static* does not scale, because the number of critical nodes is a constant function of the application data flow graph, and is not proportional to the size of the field. The slight dip in performance as field size increases, affecting all the algorithms, is a diminishing edge effect: events only occur within the field, and are therefore not equally distributed around nodes near the edges, making edge nodes less likely candidates for running components.

Overall, *NetPull* and *NetCenter* achieve improved system longevity by converting more of the total energy deployed in the field into useful work. The *Static* approach suffers from hotspots around critical nodes, severely limiting its lifetime. *Random* object placement fares better than *Static* because it distributes the load more evenly across the network. However, since it is not informed by application behavior, it expends unnecessary energy in each epoch, limiting the total system lifetime.

5.6 Validation of Code Rewriting Step

The simulation results above show how *NetCenter* and *NetPull* operate over large networks. In this section we evaluate the overhead of code rewriting using the real system. While the automatic code partitioning performed by our system can be overridden with a manual partitioning where necessary, we would like to characterize the performance impact of automatically rewriting an application without any programmer involvement.

For this benchmark, we ran both a hand-coded and an automatically converted version of the same application. We used two laptops with 200MHz Pentium 2 processors and 128Mb of RAM, running RedHat Linux 7.0 (2.4.2 kernel) and the Sun Java 2 Hotspot Client VM 1.3.1. We ran our application over a wired 100Mbps LAN with a single network hop between the two machines. We ran the same experiment over an 11Mb wireless 802.11b LAN. While the wireless LAN is the more realistic scenario, we have also included the wired LAN numbers, because the wired LAN setup has faster round-trip times, which would emphasize any overhead in performance.

Benchmark	Wired LAN	Wireless LAN
Automatically generated	28ms \pm 15	29ms \pm 17
Manually generated	27ms \pm 17	28ms \pm 16

Figure 10: Automatically generated code runs comparably with manually generated code.

Figure 10 shows the average of thirty execution times of our benchmark application. Even though traditional application performance is of secondary importance compared to power conservation in the context of ad hoc and sensor networks, we note that the performance of our automatically generated code is comparable to that of the manually generated version written by an experienced programmer. The large variance in execution time is primarily due to thread scheduling and synchronization of the user-level network stack.

6 Future Work

In the near future, we plan to extend this work in three directions. First, we have focused on highly dynamic versions of *NetCenter* and *NetPull*, which aggressively move components as soon as they identify an opportunity to save power. In some applications, where the communication patterns form a time-varying graph, such an aggressive approach can lead to inefficient, cyclic behavior. We intend to study variants of these algorithms with differing amounts of hysteresis to dampen feedback loops, and with additional input parameters. Second, we plan to extend our simulations to include node mobility.

While many terrestrial ad hoc sensor networks will remain static once deployed, some could exhibit mobility. Based on our current results, which show that these algorithms respond well to frequent node failure, we believe that *NetCenter* and *NetPull* will adapt well to changes in network topology. We plan to add node mobility to our simulator and examine the results firsthand. Finally, we plan to investigate the use of our single system image Java operating system with more diverse access patterns and memory allocation. We currently perform local garbage collection on each node, but do not perform distributed garbage collection. We intend to integrate a standard distributed garbage collection algorithm into our system.

7 Conclusion

In this paper, we introduced a new operating system designed for ad hoc and sensor networks. Our system enables power-aware, adaptive and easy-to-develop ad hoc networking applications by providing a single Java virtual machine image of the network. We provide an application partitioning service that automatically and efficiently converts monolithic applications into distributed components. Our system then transparently manages the placement of these application components in the network to efficiently manage the available resources.

We have also introduced *NetPull* and *NetCenter*, two power-aware, online algorithms used within our system for automatic object placement within an ad hoc network. They are practical, entail low overhead and are easy to implement because they rely only on local information that is easily available from the network. Both algorithms shorten the mean path length of data packets by observing application communication behavior and moving communicating objects closer together. This not only saves energy, but also avoids communication hotspots that lead to early application failure. Simulation results show that these automated methods for dynamic object placement can achieve a factor of four to five improvement in system longevity over static placement techniques, while actual measurements from the system demonstrate that they entail little runtime overhead.

In general, the ad hoc networking domain is rapidly emerging, with few established mechanisms, policies, benchmarks or even applications as yet. We believe that high-level abstractions, such as single system image operating systems combined with automatic object migration algorithms, will create a efficient and familiar programming environment, thereby enabling rapid development of platform-independent, adaptive ad hoc applications.

8 References

- [Amiri 00] Khalil Amiri, David Petrou, Greg Ganager and Garth Gibson. Dynamic Function Placement in Active Storage Clusters. In *Proceedings of the USENIX Annual Technical Conference*, San Diego, CA, June 2000
- [Aridor et al. 99] Yariv Aridor, Michael Factor and Avi Teperman. cJVM: a Single System Image of a JVM on a Cluster. *IEEE International Conference on Parallel Processing*, September 1999.
- [Awerbuch 87] Baruch Awerbuch. Optimal distributed algorithms for minimum weight spanning tree, counting, leader election and related problems. In *Proceedings of the 19th ACM Symposium on Theory of Computing*, 230-240, 1987.
- [Birrell & Nelson 84] A. D. Birrell and B. J. Nelson. Implementing remote procedure calls. *ACM Transactions on Computer Systems*, 2(1):39--59, February 1984.
- [Broch et al. 96] J. Broch, D. B. Johnson, and D. A. Maltz, The Dynamic Source Routing Protocol for Mobile Ad hoc Networks. *Mobile Computing*, 353, 1996.
- [Brunekreef et al. 96] J. Brunekreef, J.-P. Katoen, R. Koymans and S. Mauw, Design and analysis of dynamic leader election protocols in broadcast networks. *Distributed Computing*, 9 (4) 157-171, 1996
- [Chang and Roberts 79] E.G. Chang and R. Roberts. An improved algorithm for decentralized extremafinding in circular configuration of processors. *Communications of the ACM*. 22 (5) 281-3, 1979.
- [Cheriton 88] David Cheriton. The V Distributed System. *Communications of the ACM*, 31(3), March 1988, pp.314-333.
- [Garcia-Molina 82] Hector Garcia-Molina. Elections in Distributed Computer Systems. *IEEE Transactions on Computers*. C-31 (1) 48-59, 1982.
- [Gehani 97] Samir B. Gehani. xDU: A Java-based Framework for Explicitly Partitioning Applications into Distributable Units. Master of Science Thesis, Department of Computer Science, University of San Francisco, 1997
- [Gomez et al. 01] J. Gomez, A. T. Campbell, M. Naghshineh and C. Bisdikian. PARO: Conserving Transmission Power in Wireless Ad hoc Networks. In *Proceedings of the 9th International Conference on Network Protocols*, Riverside, California, November 2001.

- [Grunwald et al. 00] Dirk Grunwald, Philip Levis, Keith I. Farkas, Charles B. Morrey III and Michael Neufeld. Policies for Dynamic Clock Scheduling. In *Proceedings of the Fourth Symposium on Operating Systems Design and Implementation*, San Diego, California, October 2000.
- [Haas & Pearlman 98] Z. J. Haas and M. R. Pearlman, The zone routing protocol (ZRP) for ad hoc networks (Internet-Draft). Mobile Ad hoc Network (MANET) Working Group, IETF, Aug. 1998.
- [Harold 00] Harold, E. R. Java Network Programming. O'Reilly & Associates, August 2000.
- [Heidemann et al. 01] John Heidemann, Fabio Silva, Chalmarek Intanagonwiwat, Ramesh Govindan, Deborah Estrin and Deepak Ganesan. Building efficient wireless sensor networks with low-level naming. In *Proceedings of the 18th Symposium on Operating Systems Principles*, Lake Louise, Alberta, October 2001.
- [Hunt & Scott 99] Galen C. Hunt and Michael L. Scott. The Coign Automatic Distributed Partitioning System. In *Proceedings of the Third Symposium on Operating System Design and Implementation*, pp. 187-200. New Orleans, Louisiana, February 1999.
- [Jacquet et al. 98] P. Jacquet, P. Muhlethaler and A. Qayyum. Optimized Link State Routing Protocol. Internet Draft, draftietf-manet-olsr-00.txt, 1998.
- [Joseph et al. 95] Anthony D. Joseph, Alan F. De Lespinasse, Joshua A. Tauber, David K. Gifford, and M. Frans Kaashoek., Rover: A Toolkit for Mobile Information Access. In *Proceedings of the Fifteenth Symposium on Operating System Principles*, December 1995.
- [Jul et al. 88] Eric Jul, Henry Levy, Norman Hutchinson, Andrew Black. Fine-Grained Mobility in the Emerald System. *ACM Transactions on Computer Systems*, 6(1), Feb. 1988, 109-133.
- [Ko & Vaidya 00] Y.-B. Ko and N. H. Vaidya, GeoTORA: A protocol for geocasting in mobile ad hoc networks. Tech. Rep. 00-010, Dept. of Computer Science, Texas A&M University, March 2000.
- [Lewis & Grimshaw 95] Mike Lewis and Andrew Grimshaw, The Core Legion Object Model, *Proceedings of the Fifth IEEE International Symposium on High Performance Distributed Computing*, August 1995

- [Litzkow et al. 97] Michael Litzkow, Todd Tannenbaum, Jim Basney, and Miron Livny. Checkpoint and migration of UNIX processes in the Condor distributed processing system. Technical Report Computer Sciences Technical Report #1346, University of Wisconsin-Madison, April 1997.
- [MANET 01] IETF Mobile Ad hoc Networking Working Group. Mobile Ad hoc Networking (MANet). http://tonnant.itd.nrl.navy.mil/manet/manet_home.html.
- [Ma et al. 99] Matchy J. M. Ma, Cho-Li Wang, Francis C. M. Lau and Zhiwei Xu. JESSICA: Java-Enabled Single System Image Computing Architecture. *The International Conference on Parallel and Distributed Processing Techniques and Applications*, June 1999.
- [Malpani et al. 00] Navneet Malpani, Jennifer L. Welch and Nitin Vaidya. Leader Election Algorithms for Mobile Ad hoc Networks. In *Proceedings of the 4th International Workshop on Discrete Algorithms and Methods for Mobile Computing and Communications*. 96--103, August 2000.
- [Mascolo 01] Cecilia Mascolo, Licia Capra and Wolfgang Emmerich. XMIDDLE - A Middleware of Ad hoc Networks. UCL-CS Research Note 00/54, 2001.
- [McCanne & Floyd 01] Steven McCanne and Sally Floyd, UCB/LBNL/VINT Network Simulator - ns (version 2), <http://www.wmash.cs.berkeley.edu/ns/>.
- [Ostrovsky 94] Rafail Ostrovsky, Sridhar Rajagopalan and Umesh Vazirani. Simple and efficient leader election in the full information model. In *Proceedings of the 26th annual ACM Symposium on Theory of Computing*, 234-42, May 1994
- [Ousterhout 88] J. Ousterhout, A. Cherenon, F. Dougliis, M. Nelson, and B. Welch. The Sprite network operating system. *IEEE Computer*, 21(2):23-36, February 1988.
- [Park & Corson 98] Vincent D. Park and M. Scott Corson. Temporally-Ordered Routing Algorithm (TORA) version 1: Functional Specification. Internet-Draft, draft-ietf-manet-tora-spec01.txt, August 1998.
- [Pei et al. 00] G. Pei, M. Gerla and X. Hong. LANMAR: Landmark Routing for Large Scale Wireless Ad hoc Networks with Group Mobility. In *Proceedings of IEEE/ACM MobiHOC 2000*, Boston, MA, Aug. 2000.
- [Perkins & Bhagwat 94] C. Perkins and P. Bhagwat. Highly dynamic destination-sequenced distance-vector routing (DSDV) for mobile computers. In *Proceedings of the ACM SIGCOMM*, October 1994.
- [Perkins 97] Perkins, C.E. Ad hoc On-Demand Distance Vector (AODV) Routing. IETF MANET, Internet Draft, Dec.1997.

- [Plank et al. 95] James S. Plank, Micah Beck, Gerry Kingsley, and Kai Li. Libckpt: Transparent Checkpointing under Unix. In *Proceedings of the Usenix Winter 1995 Technical Conference*, New Orleans, LA, January, 1995.
- [Popek & Walker 85] G. Popek and B. Walker, eds. *The LOCUS Distributed System Architecture*. MIT Press, Cambridge, MA 1985.
- [Royer & Toh 99] Elizabeth Royer and C.-K. Toh. A review of current routing protocols for ad-hoc mobile wireless networks. *IEEE Personal Communications Magazine*, April 1999, 46-55.
- [Russell et al. 99] Alexander Russell, Michael Saks and David Zuckerman. Lower bounds for leader election and collective coin-flipping in the perfect information model. In *Proceedings of the 31st ACM Symposium on Theory of Computing*. 339-47, May 1999.
- [Satyanarayanan 96] M. Satyanarayanan. Fundamental Challenges in Mobile Computing. In *Proceedings of the Fifteenth ACM Symposium on Principles of Distributed Computing*. Philadelphia, PA, May 1996.
- [Sayeed et al. 95] H.M. Sayeed, M. Abu-Amara and H. Abu-Amara. Optimal asynchronous agreement and leader election algorithm for complete networks with byzantine faulty links. *Distributed Computing*, 9 (3) 147-156, 1995.
- [Singh 96] G. Singh, Leader election in the presence of link failures. *IEEE Transactions on Parallel and Distributed Systems*, 7 (3). 231-236, March 1996.
- [Sirer et al. 99] Emin Gun Sirer, Robert Grimm, Arthur J. Gregory and Brian N. Bershad. Design and Implementation of a Distributed Virtual Machine for Networked Computers. In *Proceedings of the Seventeenth Symposium on Operating Systems Principles*, pages 202-216, Kiawah Island, South Carolina, December 1999.
- [Sivakumar et. al. 99] Raghupathy Sivakumar, Prasun Sinha and Vaduvur Bharghavan. CEDAR: a core-extraction distributed ad hoc routing algorithm. *IEEE Journal on Selected Areas in Communication*. Vol 17, No. 8, August 1999.
- [Steketee 95] Chris Steketee, Piotr Socko, Bartosz Kiepuszewski. Experiences with the Implementation of a Process Migration Mechanism for Amoeba. In *Proceedings of the 19th Australasian Computer Science Conference*, January 1995, pp. 213-224.
- [Stellner 96] Georg Stellner. CoCheck: Checkpointing and Process Migration for MPI. In *Proceedings of the International Parallel Processing Symposium*, pp. 526--531, Honolulu, HI, April 1996.

- [Tennenhouse & Wetherall 96] D. L. Tennenhouse and D. Wetherall. Towards an Active Network Architecture. In *Multimedia Computing and Networking*, San Jose, California, January 1996.
- [Weiser et al. 94] Mark Weiser, Brent Welch, Alan Demers, and Scott Shenker. Scheduling for reduced CPU energy. In *Proceedings of the First Symposium on Operating Systems Design and Implementation*, pp. 13-23, Monterey, California, November 1994.
- [Wu & Tay 99] C.W. Wu and Y.C. Tay. AMRIS: A Multicast Protocol for Ad hoc Wireless Networks. In *Proceedings of the IEEE Military Communications Conference (MILCOM)*, Atlantic City, NJ, November 1999, pp. 25-29.
- [Zayas 87] E. Zayas. Breaking the Process Migration Bottleneck. In *Proceedings of the Eleventh Symposium on Operating Systems Principles*, pp. 13--22, Austin, TX, November 1987.